

# Introduction to Hidden Markov Model and Its Application

April 16, 2005

Dr. Sung-Jung Cho

[sung-jung.cho@samsung.com](mailto:sung-jung.cho@samsung.com)

Samsung Advanced Institute of Technology (SAIT)

April 16, 2005, S.-J. Cho

1

## Contents

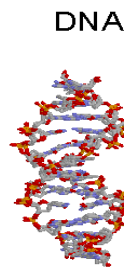
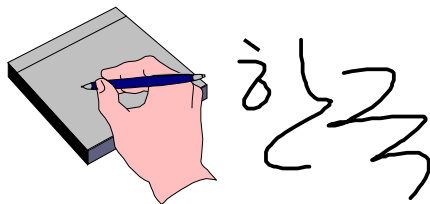
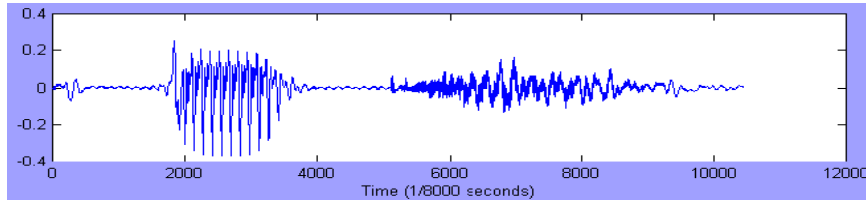
- Introduction
- Markov Model
- Hidden Markov model (HMM)
- Three algorithms of HMM
  - Model evaluation
  - Most probable path decoding
  - Model training
- Pattern classification by HMM
- Application of HMM to on-line handwriting recognition with HMM toolbox for Matlab
- Summary
- References

April 16, 2005, S.-J. Cho

2

# Sequential Data

- Data are sequentially generated according to time or index
- Spatial information along time or index



```

AAAACAAAAGCTTACAAGATGAGACATGATAAAAGCCTCCATTTC
AGGTTAGGTAATATGCTTTGGTATCCCTGTAGTTAAAAGCTTTTTC
TCTTATTTTAGAATACGTGACTATTTCTTTAGTATTAATTTTTC
CTTCTGTTTTCCATCAGGGAACCCCAAGAGCATCCAATAGAA
GCTGTGCAATTAATGTAATAATTTCAACTGCTCTTCTCAAAATAAA
GAAATATGCTAATCTTTACTCTGTATACAGTGCAGAGCCTTCTCAG
AAGCAGAAATATTTTATATTTTCCCTTATGFGAATTTTAAAGCT
GCAAACTGATGGCTTAATTTCCCTTTTGCACACTGAAAGTTTC
TAAAAGAAATCATGCTCATACACTTTGTTCCAGATGCAATTAAT
TGACACTGAACTTAATAAGTGTACTCTTTCGGAAGGCTTCTTC
AAATTTTGGACTTTTTCGATGCTTTTCTTTTCTTTTAA
ACTTCTIATGAGGAGGCGGTAATAAACCACCTGCGCTTGG
TGTAAATGAGATTGCCCATCTAGACTAGCAATCTCTCATTA
TTCCTGCTATATAAAAACCTGCTGAGGAGGCGGAAAGCA
TTTTCAATAATGCACTTTTCTACTGAAATTTTTTCTAATAAC
CAATCAAGCTTATAATTTTTTAAAATAGAAATTTCTAAGAAC
GCAATATTAACCTAATCACCATGAAAGCACTCTGGATGATGGATT
CGACAAAAGCTGGTTTTATGGTACTCTCTCTTAGATCTTAA
TTCAATGAGGAGGCTGGGCGAGGAGGCTGGAGCGGCGAAGGCTT
CTCTATTAATAATGCAATTCCTTCTTTTTTAAAGATAGCTAACTT
GCTAAATTTCTTATCTGACATTAACAATAAAAAGCTCTTTTAA
TATTAGATAA
    
```

18

April 16, 2005, S.-J. Cho

3

## Advantage of HMM on Sequential Data

- Natural model structure: doubly stochastic process
  - transition parameters model *temporal* variability
  - output distribution model *spatial* variability
- Efficient and good modeling tool for
  - sequences with temporal constraints
  - spatial variability along the sequence
  - real world complex processes
- Efficient evaluation, decoding and training algorithms
  - Mathematically strong
  - Computationally efficient
- Proven technology!
  - Successful stories in many applications

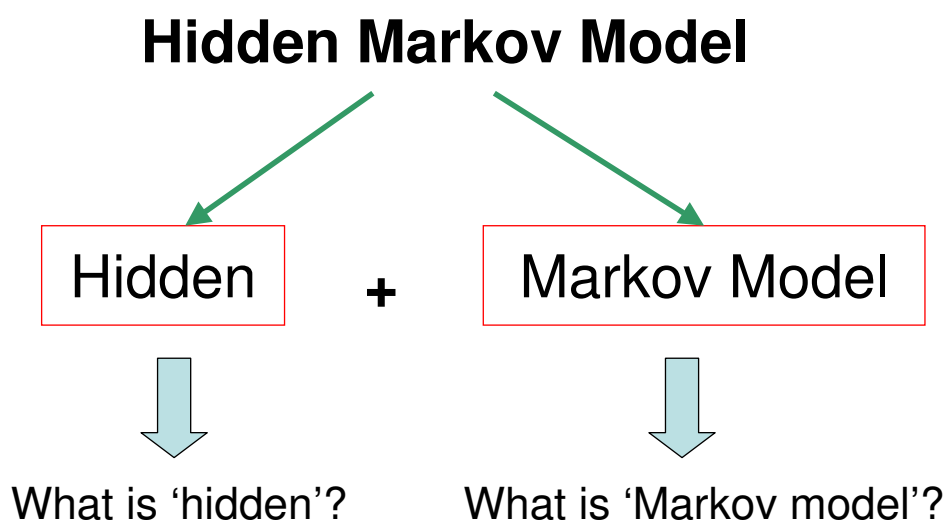
April 16, 2005, S.-J. Cho

4

## Successful Application Areas of HMM

- On-line handwriting recognition
- Speech recognition
- Gesture recognition
- Language modeling
- Motion video analysis and tracking
- Protein sequence/gene sequence alignment
- Stock price prediction
- ...

## What's HMM?



# Markov Model

- Scenario
- Graphical representation
- Definition
- Sequence probability
- State probability

April 16, 2005, S.-J. Cho

7

## Markov Model: Scenario

- Classify a weather into three states
  - State 1: rain or snow
  - State 2: cloudy
  - State 3: sunny
- By carefully examining the weather of some city for a long time, we found following weather change pattern



		Tomorrow		
		Rain/snow	Cloudy	Sunny
Today	Rain/Snow	0.4	0.3	0.3
	Cloudy	0.2	0.6	0.2
	Sunny	0.1	0.1	0.8

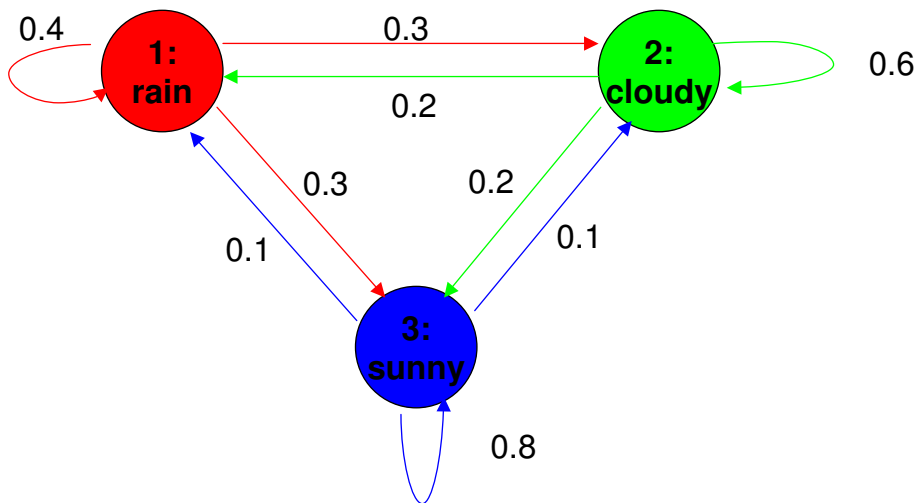
Assumption: tomorrow weather depends only on today one!

April 16, 2005, S.-J. Cho

8

# Markov Model: Graphical Representation

- Visual illustration with diagram



- Each state corresponds to one observation
- Sum of outgoing edge weights is one

April 16, 2005, S.-J. Cho

9

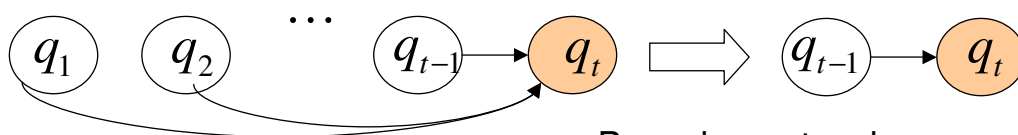
## Markov Model: Definition

- Observable states  
 $\{1, 2, \dots, N\}$
- Observed sequence

$$q_1, q_2, \dots, q_T$$

- 1<sup>st</sup> order Markov assumption

$$P(q_t = j \mid q_{t-1} = i, q_{t-2} = k, \dots) = P(q_t = j \mid q_{t-1} = i)$$



Bayesian network representation

- Stationary

$$P(q_t = j \mid q_{t-1} = i) = P(q_{t+l} = j \mid q_{t+l-1} = i)$$

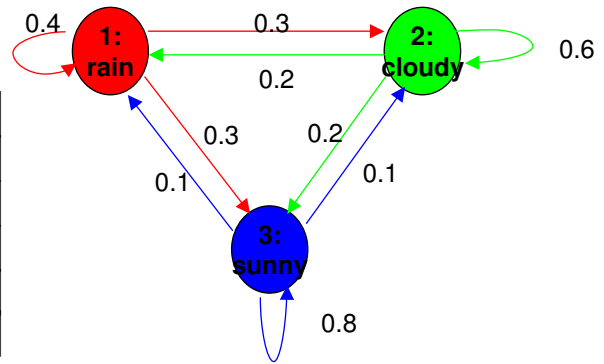
April 16, 2005, S.-J. Cho

10

## Markov Model: Definition (Cont.)

- State transition matrix

$$A = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1N} \\ a_{21} & a_{22} & \cdots & a_{2N} \\ \vdots & \vdots & \vdots & \vdots \\ a_{N1} & a_{NN} & \cdots & a_{NN} \end{bmatrix}$$



- Where

$$a_{ij} = P(q_t = j \mid q_{t-1} = i), \quad 1 \leq i, j \leq N$$

- With constraints

$$a_{ij} \geq 0, \quad \sum_{j=1}^N a_{ij} = 1$$

- Initial state probability

$$\pi_i = P(q_1 = i), \quad 1 \leq i \leq N$$

## Markov Model: Sequence Prob.

- Conditional probability

$$P(A, B) = P(A \mid B)P(B)$$

- Sequence probability of Markov model

$$\begin{aligned} & P(q_1, q_2, \dots, q_T) \\ & \xrightarrow{\text{Chain rule}} \\ & = P(q_1)P(q_2 \mid q_1) \cdots P(q_{T-1} \mid q_1, \dots, q_{T-2})P(q_T \mid q_1, \dots, q_{T-1}) \\ & = P(q_1)P(q_2 \mid q_1) \cdots P(q_{T-1} \mid q_{T-2})P(q_T \mid q_{T-1}) \end{aligned}$$



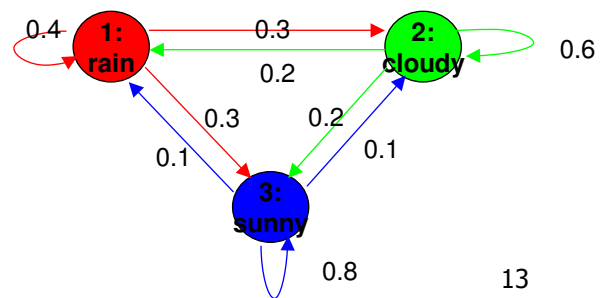
1<sup>st</sup> order Markov assumption

## Markov Model: Sequence Prob. (Cont.)

- Question: What is the probability that the weather for the next 7 days will be “sun-sun-rain-rain-sun-cloudy-sun” when today is sunny?

$S_1 : \text{rain}, S_2 : \text{cloudy}, S_3 : \text{sunny}$

$$\begin{aligned}
 P(O \mid \text{model}) &= P(S_3, S_3, S_3, S_1, S_1, S_3, S_2, S_3 \mid \text{model}) \\
 &= P(S_3) \cdot P(S_3 \mid S_3) \cdot P(S_3 \mid S_3) \cdot P(S_1 \mid S_3) \\
 &\quad \cdot P(S_1 \mid S_1) P(S_3 \mid S_1) P(S_2 \mid S_3) P(S_3 \mid S_2) \\
 &= \pi_3 \cdot a_{33} \cdot a_{33} \cdot a_{31} \cdot a_{11} \cdot a_{13} \cdot a_{32} \cdot a_{23} \\
 &= 1 \cdot (0.8)(0.8)(0.1)(0.4)(0.3)(0.1)(0.2) \\
 &= 1.536 \times 10^{-4}
 \end{aligned}$$

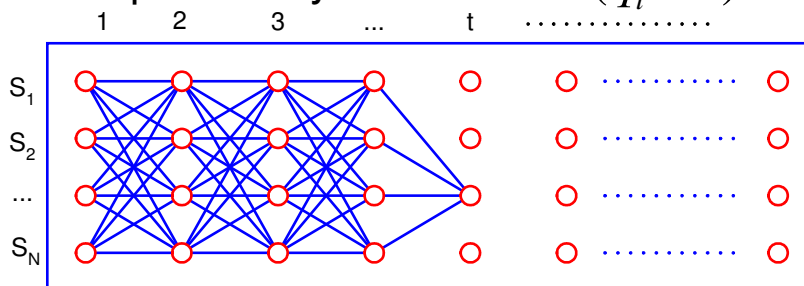


April 16, 2005, S.-J. Cho

13

## Markov Model: State Probability

- State probability at time  $t$ :  $P(q_t = i)$



- Simple but slow algorithm:

– Probability of a path that ends to state  $i$  at time  $t$ :

$$Q_t(i) = (q_1, q_2, \dots, q_t = i)$$

$$P(Q_t(i)) = \pi_{q_1} \prod_{k=2}^t P(q_k \mid q_{k-1})$$

– Summation of probabilities of all the paths that ends to  $i$  at  $t$

$$P(q_t = i) = \sum_{\text{all } Q_t(i)\text{'s}} P(Q_t(i))$$

Exponential time complexity:

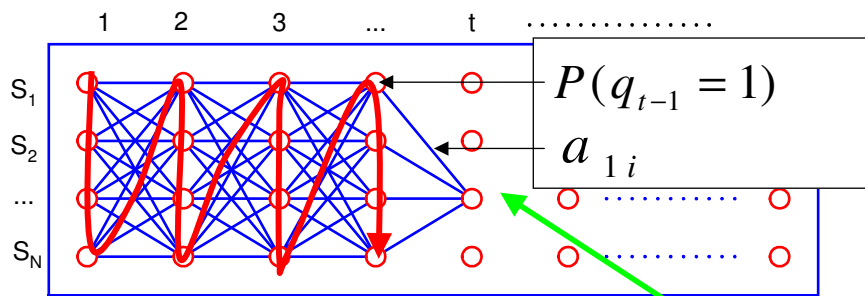
$$O(N^t)$$

April 16, 2005, S.-J. Cho

14

## Markov Model: State Prob. (Cont.)

- State probability at time  $t$  :  $P(q_t = i)$



- Efficient algorithm

- Recursive path probability calculation

$$P(q_t = i) = \sum_{j=1}^N P(q_{t-1} = j, q_t = i)$$

$$= \sum_{j=1}^N P(q_{t-1} = j) P(q_t = i \mid q_{t-1} = j)$$

$$= \sum_{j=1}^N P(q_{t-1} = j) \cdot a_{ji}$$

April 16, 2005, S.-J. Cho

Time complexity:  $O(N^2t)$

## What's HMM?

### Hidden Markov Model

Hidden

+

Markov Model

What is 'hidden'?

What is 'Markov model'?



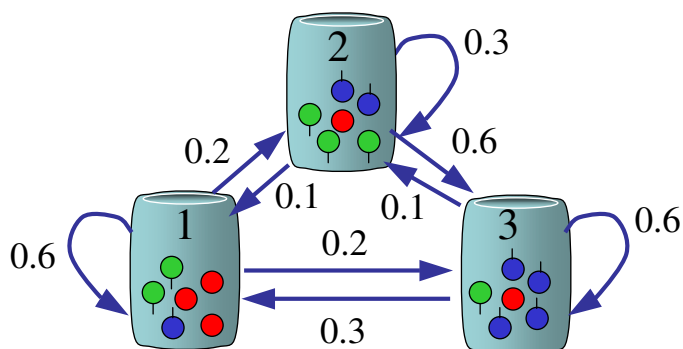
## Hidden Markov Model

- Example
- Generation process
- Definition
- Model evaluation algorithm
- Path decoding algorithm
- Training algorithm

April 16, 2005, S.-J. Cho

17

## Hidden Markov Model: Example



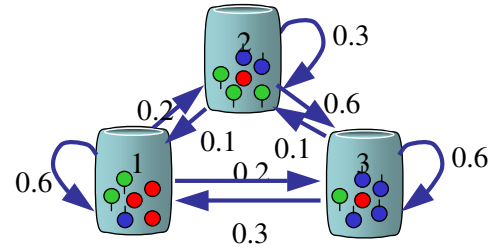
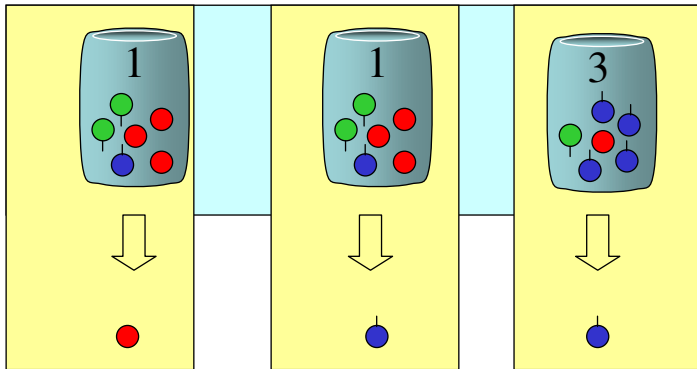
- N urns containing color balls
- M distinct colors
- Each urn contains different number of color balls

April 16, 2005, S.-J. Cho

18

# HMM: Generation Process

- Sequence generating algorithm
  - Step 1: Pick initial urn according to some random process
  - Step 2: Randomly pick a ball from the urn and then replace it
  - Step 3: Select another urn according to a random selection process
  - Step 4: Repeat steps 2 and 3

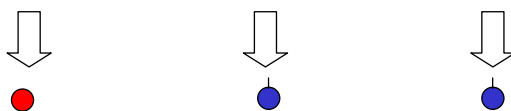
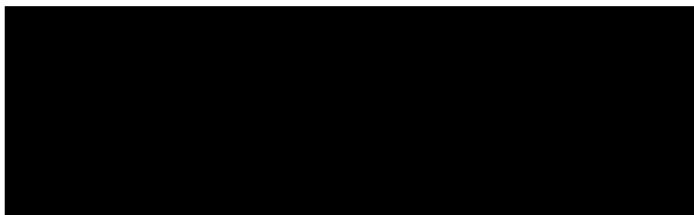


Markov process:  $\{q(t)\}$

Output process:  $\{f(x|q)\}$

# HMM: Hidden Information

- Now, what is hidden?



- We can just see the chosen balls
- We can't see which urn is selected at a time
- So, urn selection (state transition) information is hidden

## HMM: Definition

- Notation:  $\lambda = (A, B, \Pi)$

(1) N: Number of states

(2) M: Number of symbols observable in states

$$V = \{v_1, \dots, v_M\}$$

(3) A: State transition probability distribution

$$A = \{a_{ij}\}, \quad 1 \leq i, j \leq N$$

(4) B: Observation symbol probability distribution

$$B = \{b_i(v_k)\}, \quad 1 \leq i \leq N, 1 \leq k \leq M$$

(5)  $\Pi$ : Initial state distribution

$$\pi_i = P(q_1 = i), \quad 1 \leq i \leq N$$

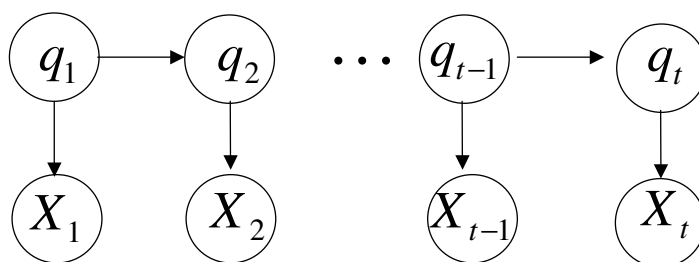
## HMM: Dependency Structure

- 1-st order Markov assumption of transition

$$P(q_t | q_1, q_2, \dots, q_{t-1}) = P(q_t | q_{t-1})$$

- Conditional independency of observation parameters

$$P(X_t | q_t, X_1, \dots, X_{t-1}, q_1, \dots, q_{t-1}) = P(X_t | q_t)$$



Bayesian network representation

# HMM: Example Revisited

- # of states:  $N=3$
- # of observation:  $M=3$   
 $V = \{R, G, B\}$

- Initial state distribution

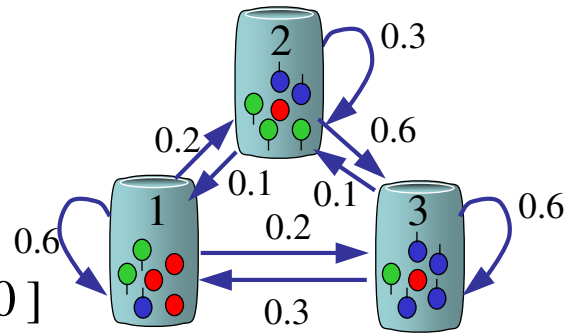
$$\pi = \{P(q_1 = i)\} = [1, 0, 0]$$

- State transition probability distribution

$$A = \{a_{ij}\} = \begin{bmatrix} 0.6 & 0.2 & 0.2 \\ 0.1 & 0.3 & 0.6 \\ 0.3 & 0.1 & 0.6 \end{bmatrix}$$

- Observation symbol probability distribution

$$B = \{b_i(v_k)\} = \begin{bmatrix} 3/6 & 2/6 & 1/6 \\ 1/6 & 3/6 & 2/6 \\ 1/6 & 1/6 & 4/6 \end{bmatrix}$$



# HMM: Three Problems

- What is the probability of generating an observation sequence?  
 – Model evaluation

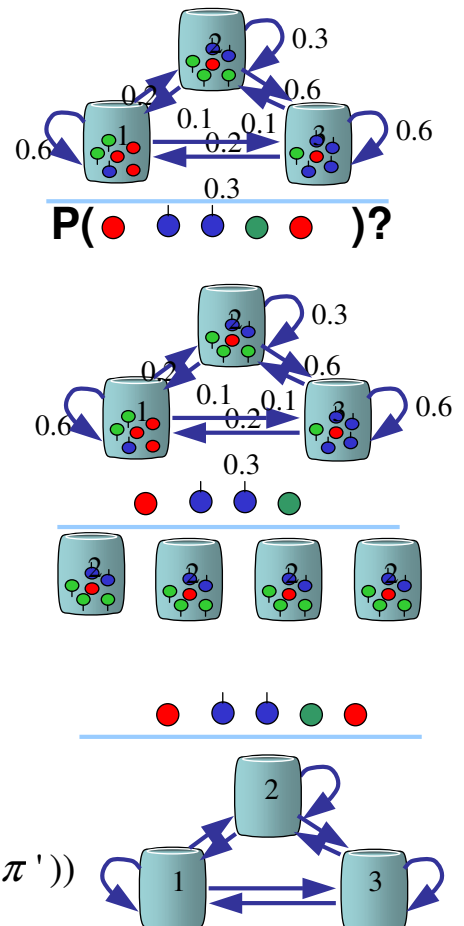
$$P(X = x_1, x_2, \dots, x_T | \lambda) = ?$$

- Given observation, what is the most probable transition sequence?  
 – Segmentation or path analysis

$$Q^* = \arg \max_{Q=(q_1, \dots, q_T)} P(Q, X | \lambda)$$

- How do we estimate or optimize the parameters of an HMM?  
 – Training problem

$$P(X | \lambda = (A, B, \pi)) < P(X | \lambda' = (A', B', \pi'))$$



# Model Evaluation

Forward algorithm  
Backward algorithm

April 16, 2005, S.-J. Cho

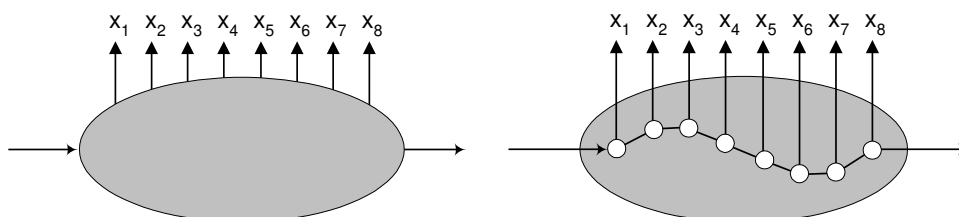
25

## Definition

- Given a model  $\lambda$
- Observation sequence:  $X = x_1, x_2, \dots, x_T$
- $P(X | \lambda) = ?$

$$P(X | \lambda) = \sum_Q P(X, Q | \lambda) = \sum_Q P(X | Q, \lambda) P(Q | \lambda)$$

(A path or state sequence:  $Q = q_1, \dots, q_T$  )



April 16, 2005, S.-J. Cho

26

## Solution

- Easy but slow solution: exhaustive enumeration

$$\begin{aligned}
 P(X | \lambda) &= \sum_Q P(X, Q | \lambda) = \sum_Q P(X | Q, \lambda) P(Q | \lambda) \\
 &= \sum_Q b_{q_1}(x_1) b_{q_2}(x_2) \cdots b_{q_T}(x_T) \pi_{q_1} a_{q_1 q_2} a_{q_2 q_3} \cdots a_{q_{T-1} q_T}
 \end{aligned}$$

- Exhaustive enumeration = combinational explosion!

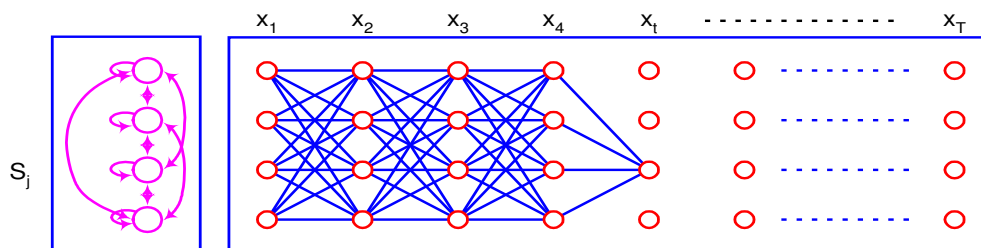
$$O(N^T)$$

- Smart solution exists?

- Yes!
- Dynamic Programming technique
- Lattice structure based computation
- Highly efficient -- linear in frame length

## Forward Algorithm

- Key idea
  - Span a lattice of N states and T times
  - Keep the sum of probabilities of all the paths coming to each state i at time t



- Forward probability

$$\begin{aligned}
 \alpha_t(j) &= P(x_1 x_2 \dots x_t, q_t = S_j | \lambda) \\
 &= \sum_{Q_t} P(x_1 x_2 \dots x_t, Q_t = q_1 \dots q_t | \lambda) \\
 &= \sum_{i=1}^N \alpha_{t-1}(i) a_{ij} b_j(x_t)
 \end{aligned}$$

# Forward Algorithm

- Initialization

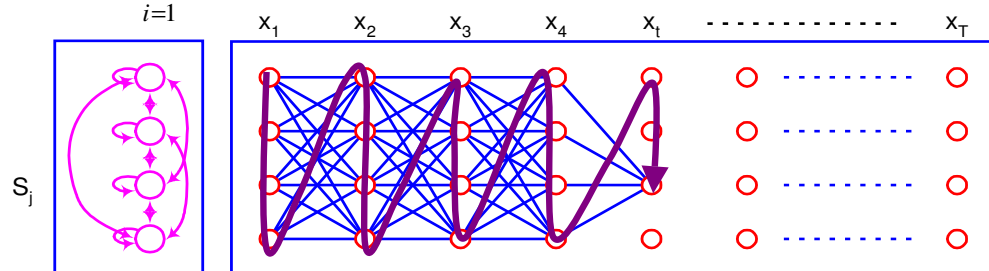
$$\alpha_1(i) = \pi_i b_i(\mathbf{x}_1) \quad 1 \leq i \leq N$$

- Induction

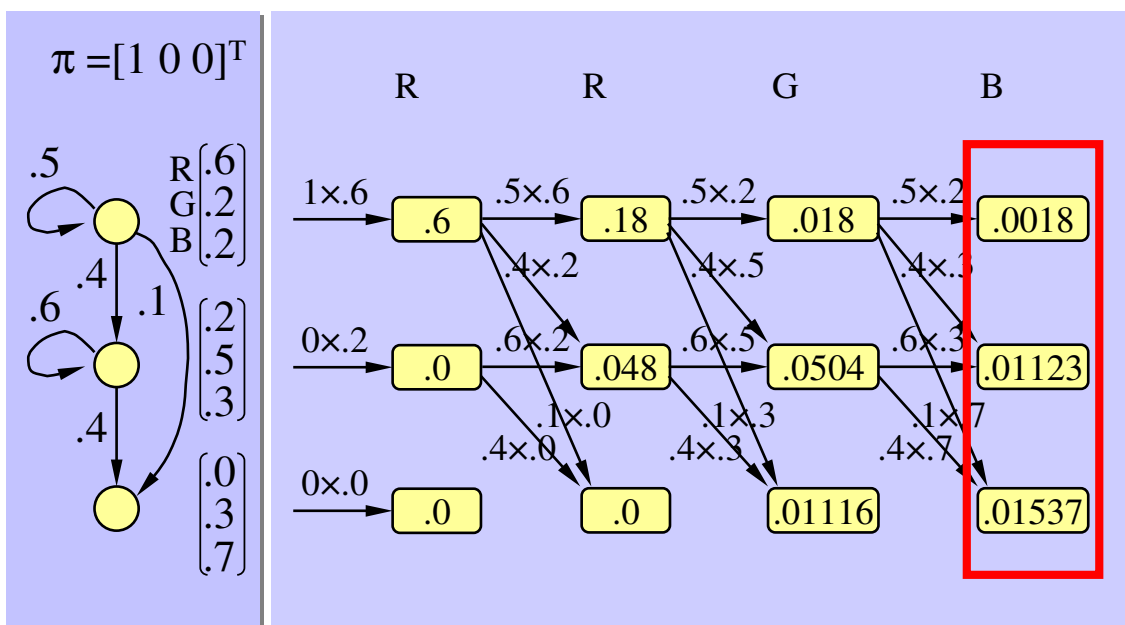
$$\alpha_t(j) = \sum_{i=1}^N \alpha_{t-1}(i) a_{ij} b_j(\mathbf{x}_t) \quad 1 \leq j \leq N, t = 2, 3, \dots, T$$

- Termination

$$P(\mathbf{X} | \lambda) = \sum_{i=1}^N \alpha_T(i)$$

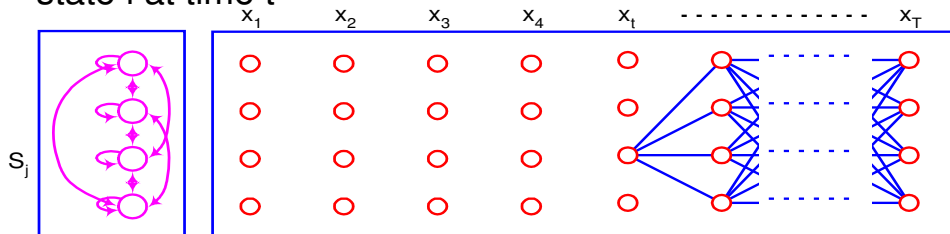


## Numerical Example: $P(\text{RRGB} | \lambda)$ [신봉기 03]



# Backward Algorithm (1)

- Key Idea
  - Span a lattice of N states and T times
  - Keep the sum of probabilities of all the outgoing paths at each state i at time t



- Backward probability

$$\begin{aligned} \beta_t(i) &= P(x_{t+1}x_{t+2}\dots x_T \mid q_t = S_i, \lambda) \\ &= \sum_{Q_{t+1}} P(x_{t+1}x_{t+2}\dots x_T, Q_{t+1} = q_{t+1}\dots q_T \mid q_t = S_i, \lambda) \\ &= \sum_{j=1}^N a_{ij} b_j(x_{t+1}) \beta_{t+1}(j) \end{aligned}$$

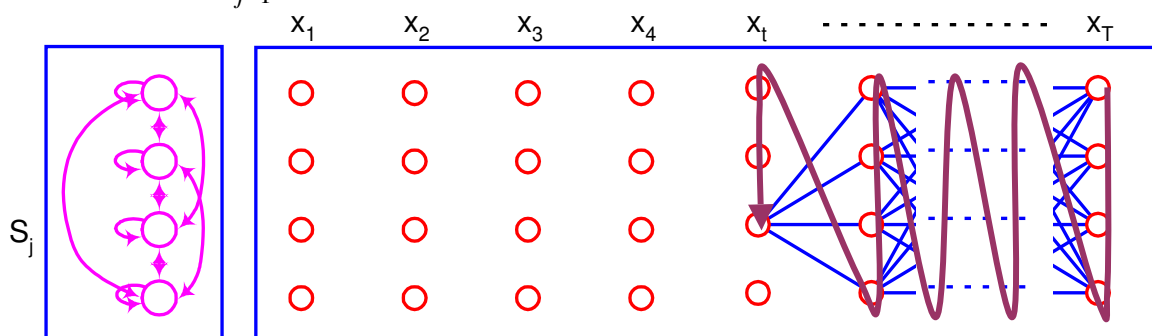
# Backward Algorithm (2)

- Initialization

$$\beta_T(i) = 1 \quad 1 \leq i \leq N$$

- Induction

$$\beta_t(i) = \sum_{j=1}^N a_{ij} b_j(\mathbf{x}_{t+1}) \beta_{t+1}(j) \quad 1 \leq i \leq N, \quad t = T-1, T-2, \dots, 1$$





# The Most Probable Path Decoding

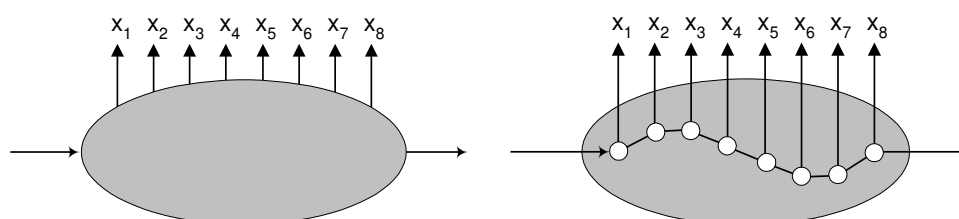
State sequence  
Optimal path  
Viterbi algorithm  
Sequence segmentation

April 16, 2005, S.-J. Cho

33

## The Most Probable Path

- Given a model  $\lambda$
- Observation sequence:  $X = \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T$
- $P(X, Q | \lambda) = ?$
- $Q^* = \arg \max_Q P(X, Q | \lambda) = \arg \max_Q P(X | Q, \lambda) P(Q | \lambda)$ 
  - (A path or state sequence:  $Q = q_1, \dots, q_T$ )



April 16, 2005, S.-J. Cho

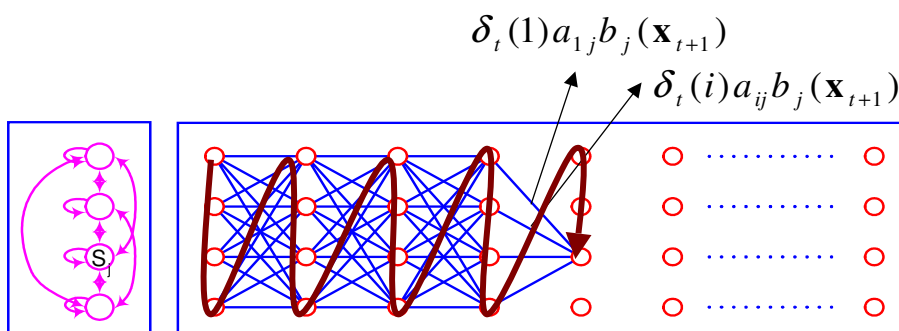
34

# Viterbi Algorithm

- Purpose
  - An analysis for internal processing result
  - The best, the most likely state sequence
  - Internal segmentation
- Viterbi Algorithm
  - Alignment of observation and state transition
  - Dynamic programming technique

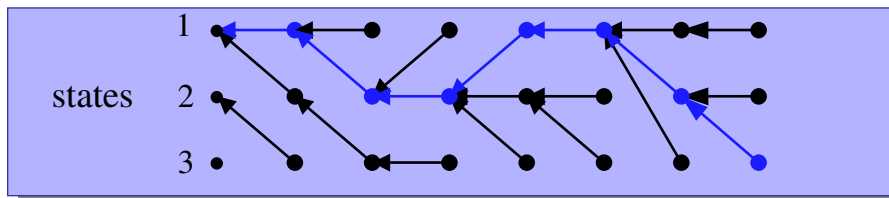
## Viterbi Path Idea

- Key idea
  - Span a lattice of N states and T times
  - Keep the probability and the previous node of the most probable path coming to each state i at time t
- Recursive path selection
  - Path probability:  $\delta_{t+1}(j) = \max_{1 \leq i \leq N} \delta_t(i) a_{ij} b_j(\mathbf{x}_{t+1})$
  - Path node:  $\psi_{t+1}(j) = \arg \max_{1 \leq i \leq N} \delta_t(i) a_{ij}$

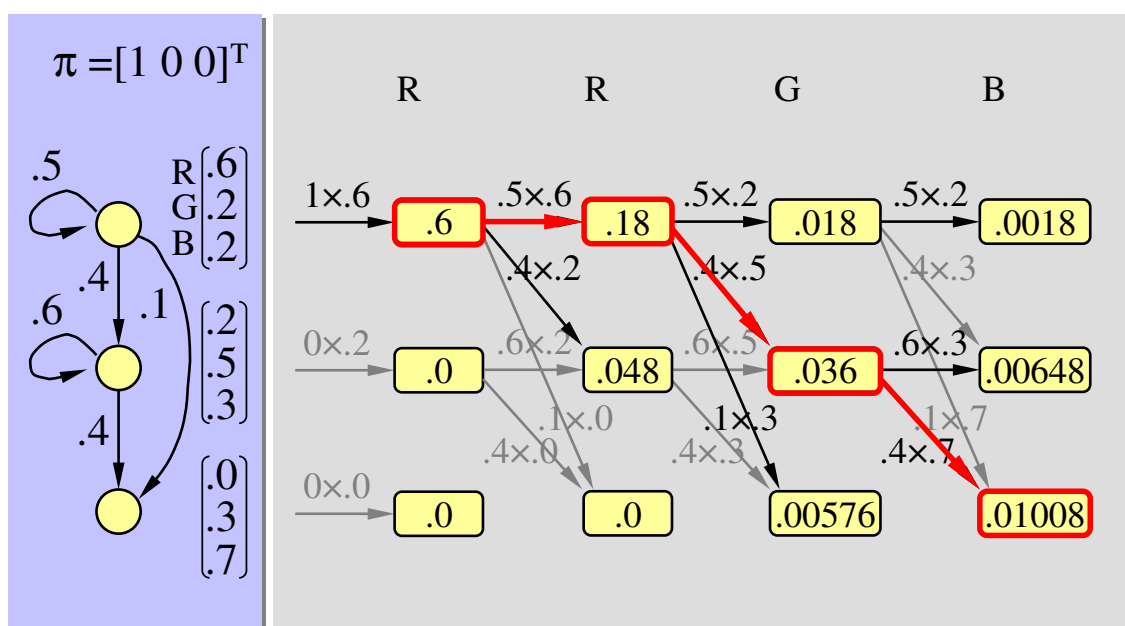


# Viterbi Algorithm

- Introduction:  $\delta_1(i) = \pi_i b_i(\mathbf{x}_1)$   
 $\psi_1(i) = 0$
- Recursion:  $\delta_{t+1}(j) = \max_{1 \leq i \leq N} \delta_t(i) a_{ij} b_j(\mathbf{x}_{t+1})$   
 $\psi_{t+1}(j) = \arg \max_{1 \leq i \leq N} \delta_t(i) a_{ij}$
- Termination:  $P^* = \max_{1 \leq i \leq N} \delta_T(i)$   
 $q_T^* = \arg \max_{1 \leq i \leq N} \delta_T(i)$
- Path backtracking:  $q_t^* = \psi_{t+1}(q_{t+1}^*), \quad t = T - 1, \dots, 1$



## Numerical Example: $P(\text{RRGB}, Q^* | \lambda)$ [신봉기 03]



# Parameter Reestimation

HMM training algorithm  
Maximum likelihood estimation  
Baum-Welch reestimation

April 16, 2005, S.-J. Cho

39

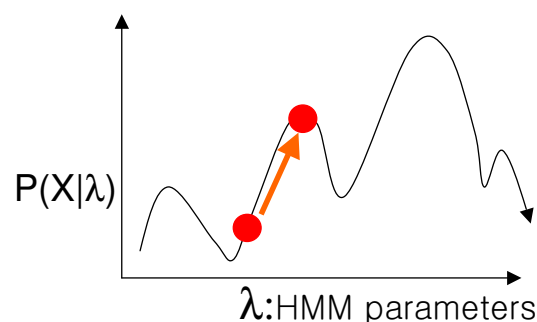
## HMM Training Algorithm

- Given an observation sequence  $X = \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T$
- Find the model parameter  $\lambda^* = (A, B, \pi)$
- s.t.  $P(X | \lambda^*) \geq P(X | \lambda)$  for  $\forall \lambda$ 
  - Adapt HMM parameters maximally to training samples
  - Likelihood of a sample

$$P(X | \lambda) = \sum_Q P(X | Q, \lambda) P(Q | \lambda)$$

State transition is hidden!

- NO analytical solution
- *Baum-Welch* reestimation (EM)
  - iterative procedures that locally maximizes  $P(X|\lambda)$
  - convergence proven
  - MLE statistic estimation



April 16, 2005, S.-J. Cho

40

# Maximum Likelihood Estimation

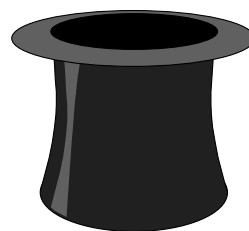
- MLE “selects those parameters that maximizes the probability function of the observed sample.”
- [Definition] Maximum Likelihood Estimate
  - $\Theta$ : a set of distribution parameters
  - Given  $X$ ,  $\Theta^*$  is maximum likelihood estimate of  $\Theta$  if
  - $f(X|\Theta^*) = \max_{\Theta} f(X|\Theta)$

April 16, 2005, S.-J. Cho

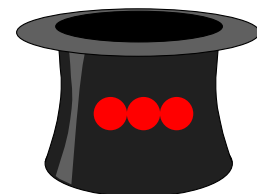
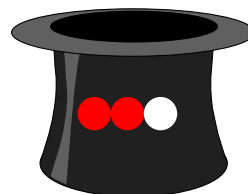
41

## MLE Example

- Scenario
  - Known: 3 balls inside urn
  - (some red; some white)
  - Unknown:  $R = \#$  red balls
  - Observation: ●● (two reds)



- Two models
  - $P(\text{●●} | R=2) = \frac{\binom{2}{2}\binom{1}{0}}{\binom{3}{2}} = \frac{1}{3}$
  - $P(\text{●●} | R=3) = \frac{\binom{3}{2}}{\binom{3}{2}} = 1$



- Which model?
  - $L(\lambda_{R=3}) > L(\lambda_{R=2})$
  - Model( $R=3$ ) is our choice

April 16, 2005, S.-J. Cho

42

## MLE Example (Cont.)

- Model( $R=3$ ) is a more likely strategy, unless we have a priori knowledge of the system.
- However, without an observation of two red balls
  - No reason to prefer  $P(\lambda_{R=3})$  to  $P(\lambda_{R=2})$
- ML method chooses the set of parameters that maximizes the likelihood of the given observation.
- It makes parameters maximally adapted to training data.

## EM Algorithm for Training

•With  $\lambda^{(t)} = \langle \{a_{ij}\}, \{b_{ik}\}, \pi_i \rangle$ , estimate EXPECTATION of following quantities:

- Expected number of state  $i$  visiting
- Expected number of transitions from  $i$  to  $j$

•With following quantities:

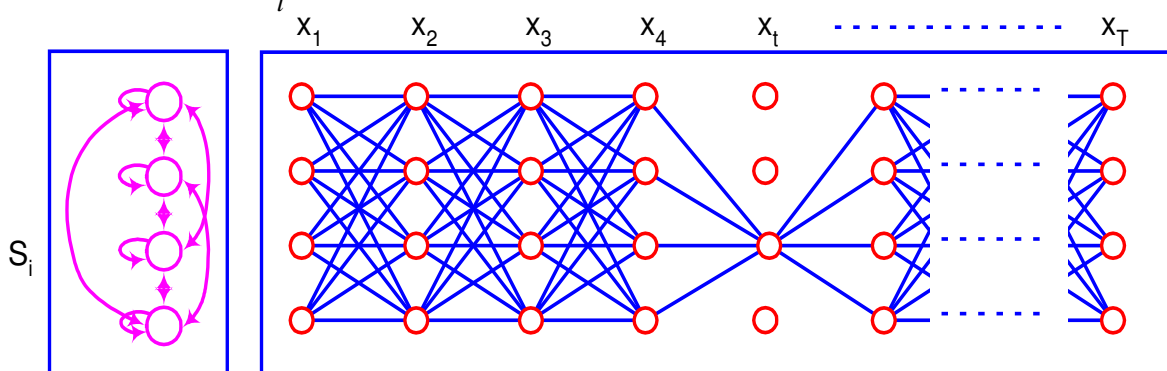
- Expected number of state  $i$  visiting
- Expected number of transitions from  $i$  to  $j$

• Obtain the MAXIMUM LIKELIHOOD of  $\lambda^{(t+1)} = \langle \{a'_{ij}\}, \{b'_{ik}\}, \pi_i \rangle$

## Expected Number of $S_i$ Visiting

$$\begin{aligned} \gamma_t(i) &= P(q_t = S_i | X, \lambda) \\ &= P(q_t = S_i, X | \lambda) / P(X | \lambda) \\ &= \frac{\alpha_t(i) \beta_t(i)}{\sum_j \alpha_t(j) \beta_t(j)} \end{aligned}$$

$$\Gamma(i) = \sum_t \gamma_t(i)$$



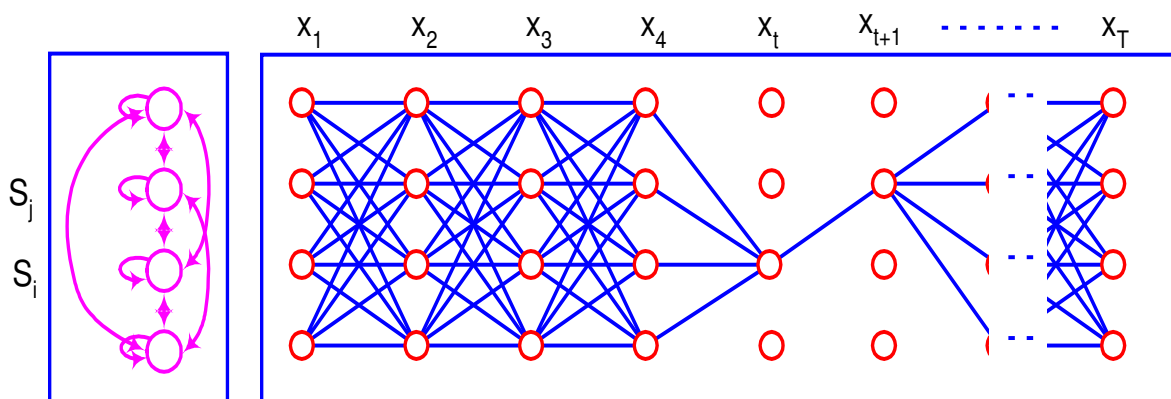
April 16, 2005, S.-J. Cho

45

## Expected Number of Transition

$$\xi_t(i, j) = P(q_t = S_i, q_{t+1} = S_j | X, \lambda) = \frac{\alpha_t(i) a_{ij} b_j(x_{t+1}) \beta_{t+1}(j)}{\sum_i \sum_j \alpha_t(i) a_{ij} b_j(x_{t+1}) \beta_{t+1}(j)}$$

$$\Xi(i, j) = \sum_t \xi_t(i, j)$$



April 16, 2005, S.-J. Cho

46

## Parameter Reestimation

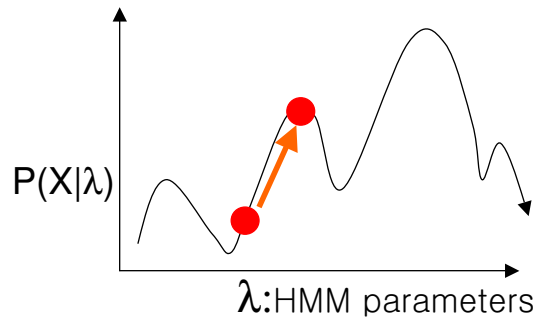
- MLE parameter estimation

$$\bar{a}_{ij} = \frac{\sum_{t=1}^{T-1} \xi_t(i, j)}{\sum_{t=1}^{T-1} \Gamma_t(i)}$$

$$\bar{b}_i(v_k) = \frac{\sum_{t=1}^T \Gamma_t(i) \delta(x_t, v_k)}{\sum_{t=1}^T \Gamma_t(i)}$$

$$\bar{\pi}_i = \gamma_1(i)$$

- Iterative:  $P(X | \lambda^{(t+1)}) \geq P(X | \lambda^{(t)})$
- convergence proven:
- arriving local optima



## Pattern Classification by HMM

- Pattern classification
- Extension of HMM structure
- Extension of HMM training method
- Practical issues of HMM
- HMM history



## Pattern Classification by HMM

- Construct one HMM per each class  $k$ 
  - $\lambda_1, \dots, \lambda_N$
- Train each HMM  $\lambda_k$  with samples  $D_k$ 
  - Baum-Welch reestimation algorithm
- Calculate model likelihood of  $\lambda_1, \dots, \lambda_N$  with observation  $X$ 
  - Forward algorithm:  $P(X | \lambda_k)$
- Find the model with maximum *a posteriori* probability

$$\begin{aligned}\lambda^* &= \operatorname{argmax}_{\lambda_k} P(\lambda_k | X) \\ &= \operatorname{argmax}_{\lambda_k} P(\lambda_k)P(X | \lambda_k) / P(X) \\ &= \operatorname{argmax}_{\lambda_k} P(\lambda_k)P(X | \lambda_k)\end{aligned}$$

## Extension of HMM Structure

- Extension of state transition parameters
  - Duration modeling HMM
    - More accurate temporal behavior
  - Transition-output HMM
    - HMM output functions are attached to transitions rather than states
- Extension of observation parameter
  - Segmental HMM
    - More accurate modeling of trajectories at each state, but more computational cost
  - Continuous density HMM (CHMM)
    - Output distribution is modeled with mixture of Gaussian
  - Semi-continuous HMM
    - Mix of continuous HMM and discrete HMM by sharing Gaussian components

## Extension of HMM Training Method

- Maximum Likelihood Estimation (MLE)\*
  - maximize the probability of the observed samples
- Maximum Mutual Information (MMI) Method
  - information-theoretic measure
  - maximize average mutual information:
$$I^* = \max_{\lambda} \left\{ \sum_{v=1}^V \left[ \log P(X^v | \lambda_v) - \log \sum_{w=1}^V P(X^w | \lambda_w) \right] \right\}$$
  - maximize discrimination power by training models together
- Minimal Discriminant Information (MDI) Method
  - minimize the DI or the cross entropy between pd(signal) and pd(HMM)'s
  - use generalized Baum algorithm

## Practical Issues of HMM

- Architectural and behavioral choices
  - the unit of modeling -- design choice
  - type of models: ergodic, left-right, etc.
  - number of states
  - observation symbols; discrete, continuous; mixture number
- Initial estimates
  - $A, \pi$  : adequate with random or uniform initial values
  - $B$  : good initial estimates are essential for CHMM

## Practical Issues of HMM (Cont.)

- Scaling

$$\alpha_t(i) = \prod_{s=1}^{t-1} a_{s,s+1} \prod_{s=1}^t b_s(x_s)$$

- heads exponentially to zero: --> scale by  $1 / \sum_{i=1}^N \alpha_t(i)$
- Multiple observation sequences
  - accumulate the expected freq. with weight  $P(X(k)|I)$
- Insufficient training data
  - deleted interpolation with desired model & small model
  - output prob. smoothing (by local perturbation of symbols)
  - output probability tying between different states

## HMM History [Roweis]

- Markov('13) and Shannon ('48, '51) studied Markov chains
- Baum et. Al (BP'66, BE'67 ...) developed many theories of “probabilistic functions of Markov chains”
- Viterbi ('67) developed an efficient optimal state search algorithm
- Application to speech recognition started
  - Baker('75) at CMU
  - Jelinek's group ('75) at IBM
- Dempster, Laird & Rubin ('77) recognized a general form of the Baum-Welch algorithm

# Application of HMM to on-line handwriting recognition with HMM SW

- SW tools for HMM
- Introduction to on-line handwriting recognition
- Data preparation
- Training & testing

## SW Tools for HMM

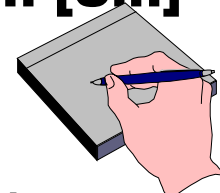
- **HMM toolbox for Matlab**
  - Developed by Kevin Murphy
  - Freely downloadable SW written in Matlab (Hmm... Matlab is not free!)
  - Easy-to-use: flexible data structure and fast prototyping by Matlab
  - Somewhat slow performance due to Matlab
  - Download: <http://www.cs.ubc.ca/~murphyk/Software/HMM/hmm.html>
- HTK (Hidden Markov toolkit)
  - Developed by Speech Vision and Robotics Group of Cambridge University
  - Freely downloadable SW written in C
  - Useful for speech recognition research: comprehensive set of programs for training, recognizing and analyzing speech signals
  - Powerful and comprehensive, but somewhat complicate and heavy package
  - Download: <http://htk.eng.cam.ac.uk/>

# SW Tools: HMM Toolbox for Matlab

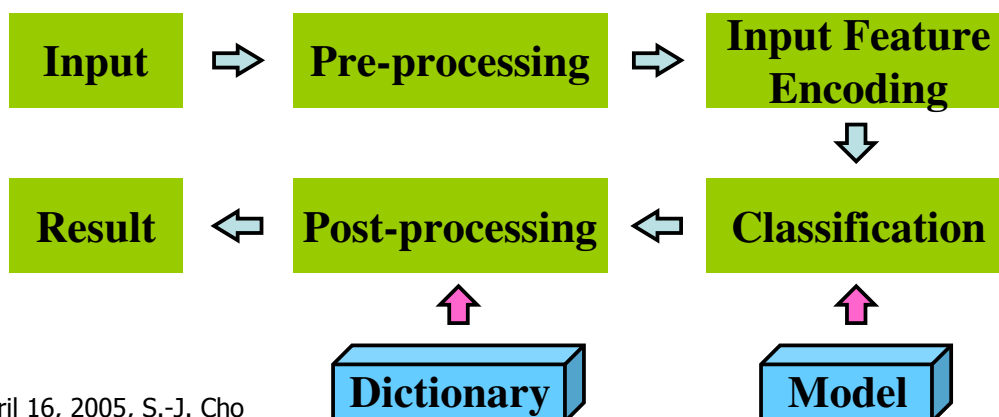
- Support training and decoding for
  - Discrete HMMs
  - Continuous HMMs with full, diagonal, or spherical covariance matrix
- 3 Algorithms for discrete HMM
  - Model evaluation (Forward algorithm)
    - $\text{Log\_likelihood} = \text{dhmm\_logprob}(\text{data}, \text{initial state probability}, \text{transition probability matrix}, \text{observation probability matrix})$
  - Viterbi decoding algorithm
    - 1)  $B = \text{multinomial\_prob}(\text{data}, \text{observation matrix});$   
 $\leftarrow B(i,t) = P(y_t | Q_t=i) \text{ for all } t,i;$
    - 2)  $[\text{path}, \text{log\_likelihood}] = \text{viterbi\_path}(\text{initial state probability}, \text{transition matrix}, B)$
  - Baum-Welch algorithm
    - $[\text{LL}, \text{prior2}, \text{transmat2}, \text{obsmat2}] = \text{dhmm\_em}(\text{data}, \text{prior1}, \text{transmat1}, \text{obsmat1}, \text{'max\_iter'}, 5);$

# On-line Handwriting Recognition [Sin]

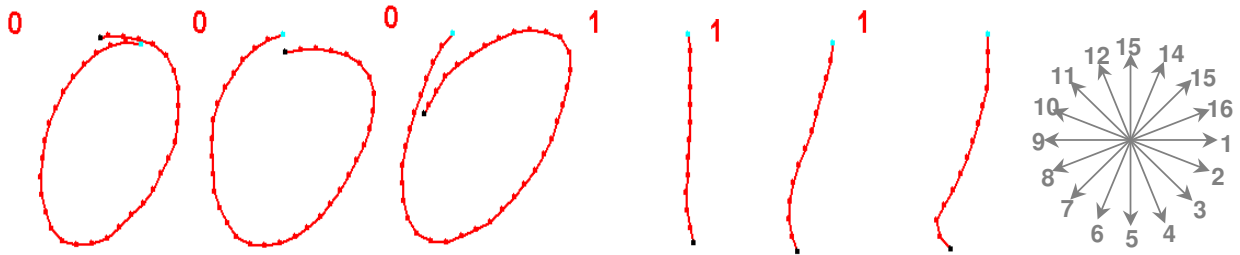
- Handwriting
  - Natural input method to human
  - Sequence of some writing units
  - Temporally ordered
    - Time series of (X,Y) ink points on tablet
- Recognition flow



한글과 나라사랑  
*the land of morning calm*  
 國民教育憲章



# Data Preparation



- Chaincode data set for class '0'
  - data0{1} = [9 8 8 7 7 7 6 6 6 5 5 5 4 4 3 2 1 16 15 15 15 15 14 14 14 13 13 13 12 12 11 10 9 9 8]
  - data0{2} = [8 8 7 7 7 6 6 5 5 5 5 4 4 3 2 1 1 16 15 15 15 15 15 14 14 14 14 13 12 11 10 10 9 9 9]
  - data0{3} = [7 6 6 6 6 6 6 5 5 5 4 3 2 1 16 16 16 15 15 15 15 14 14 14 14 14 14 13 11 10 9 9 8 8 8 7 7 6 6]
- Chaincode data set for class '1'
  - data1{1} = [5 5 5 5 5 5 5 5 5 5 4]
  - data1{2} = [5 6 6 6 6 6 6 6 6 5 5 4]
  - data1{3} = [5 5 5 6 6 6 6 6 6 7 6 4 3]

# HMM Initialization

- HMM for class '0' and randomly initialization

```

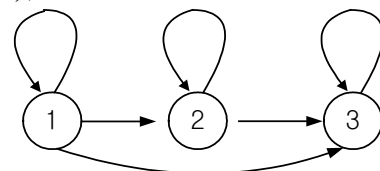
– hmm0.prior = [1 0 0];
– hmm0.transmat = rand(3,3); % 3 by 3 transition matrix
– hmm0.transmat(2,1) = 0; hmm0.transmat(3,1) = 0; hmm0.transmat(3,2) = 0;
– hmm0.transmat = mk_stochastic(hmm0.transmat);
– hmm0.transmat

```

```

0.20 0.47 0.33
0 0.45 0.55
0 0.00 1.00

```



```

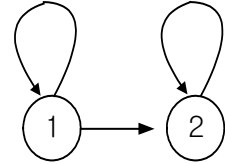
– hmm0.obsmat = rand(3, 16); % # of states * # of observation
– hmm0.obsmat = mk_stochastic(hmm0.obsmat)
0.02 0.04 0.05 0.00 0.12 0.11 0.13 0.00 0.06 0.09 0.02 0.11 0.06 0.05 0.04 0.08
0.12 0.04 0.07 0.06 0.03 0.03 0.08 0.02 0.11 0.04 0.02 0.06 0.06 0.11 0.01 0.12
0.05 0.04 0.01 0.11 0.02 0.08 0.11 0.10 0.09 0.02 0.05 0.10 0.06 0.00 0.09 0.07

```

## HMM Initialization (Cont.)

- HMM for class '1' and randomly initialization

- `hmm1.prior = [1 0];`
- `hmm1.transmat = rand(2,2); % 2 by 2 transition matrix`
- `hmm1.transmat(2,1) = 0;`
- `hmm1.transmat = mk_stochastic(hmm1.transmat);`
- `hmm1.transmat`



```

0.03 0.97
0 1.00
  
```

- `hmm1.obsmat = rand(2, 16); % # of states * # of observation`
- `hmm1.obsmat = mk_stochastic(hmm1.obsmat)`

```

0.05 0.10 0.01 0.06 0.02 0.09 0.06 0.02 0.10 0.04 0.12 0.11 0.03 0.01 0.09 0.11
0.08 0.09 0.06 0.05 0.09 0.10 0.07 0.06 0.12 0.03 0.03 0.12 0.03 0.01 0.03 0.02
  
```

## HMM Training

- Training of model 0

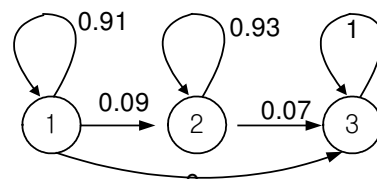
- `[LL0, hmm0.prior, hmm0.transmat, hmm0.obsmat] = dhmm_em(data0, hmm0.prior, hmm0.transmat, hmm0.obsmat)`

iteration 1, loglik = -365.390770

iteration 2, loglik = -251.112160

...

iteration 9, loglik = -210.991114



- Trained result

- `hmm0.transmat`

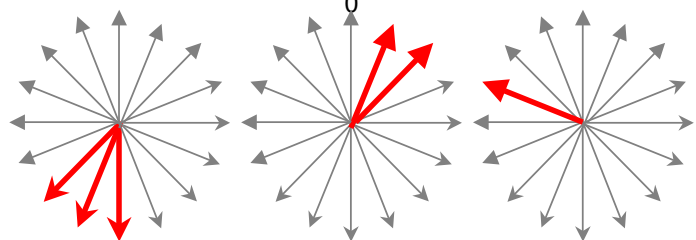
```

0.91 0.09 0.00
0.00 0.93 0.07
0.00 0.00 1.00
  
```

- `hmm0.obsmat`

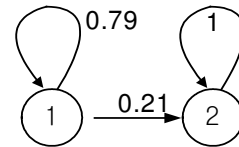
```

0.00 0.00 0.00 0.00 0.30 0.33 0.21 0.12 0.03 0.00 0.00 0.00 0.00 0.00 0.00 0.00
0.09 0.07 0.07 0.11 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.28 0.28 0.11
0.00 0.00 0.00 0.00 0.00 0.06 0.06 0.16 0.23 0.13 0.10 0.10 0.16 0.00 0.00 0.00
  
```



# HMM Training (Cont.)

- Training of model 1
  - `[LL1, hmm1.prior, hmm1.transmat, hmm1.obsmat] = dhmm_em(data1, hmm1.prior, hmm1.transmat, hmm1.obsmat)`
    - iteration 1, loglik = -95.022843
    - ...
    - iteration 10, loglik = -30.742533



- Trained model

- `hmm1.transmat`

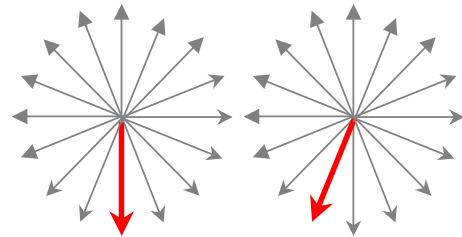
```
0.79 0.21
```

```
0.00 1.00
```

- `hmm1.obsmat`

```
0.00 0.00 0.00 0.00 1.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00
```

```
0.00 0.00 0.04 0.13 0.12 0.66 0.04 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00
```



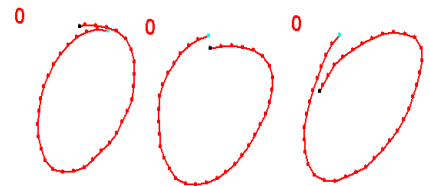
# HMM Evaluation

- Evaluation of data 0
  - for `dt = 1:length(data0)`
  - `loglike0 = dhmm_logprob(data0{dt}, hmm0.prior, hmm0.transmat, hmm0.obsmat);`
  - `loglike1 = dhmm_logprob(data0{dt}, hmm1.prior, hmm1.transmat, hmm1.obsmat);`
  - `disp(sprintf('[class 0: %d-th data] model 0: %.3f, model 1: %.3f', dt, loglike0, loglike1));`
  - end

```
[class 0: 1-th data] model 0: -68.969, model 1: -289.652
```

```
[class 0: 2-th data] model 0: -66.370, model 1: -291.671
```

```
[class 0: 3-th data] model 0: -75.649, model 1: -310.484
```

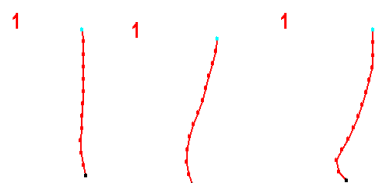


- Evaluation of data 1

```
[class 0: 1-th data] model 0: -18.676, model 1: -5.775
```

```
[class 0: 2-th data] model 0: -17.914, model 1: -11.162
```

```
[class 0: 3-th data] model 0: -21.193, model 1: -13.037
```

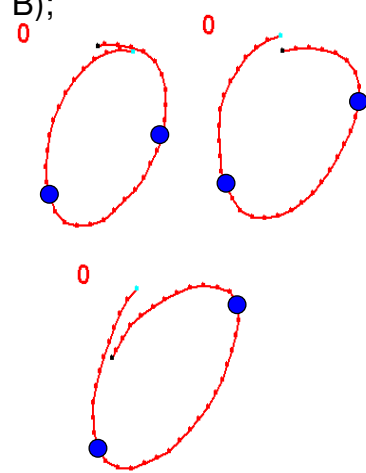




# HMM Decoding

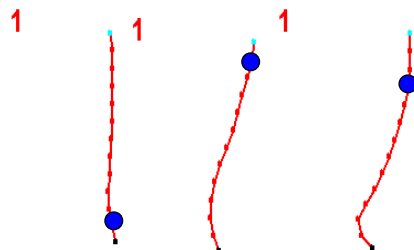
- For data '0', get the most probable path
  - for dt = 1:length(data0)
  - B = multinomial\_prob(data0{dt}, hmm0.obsmat);
  - path = viterbi\_path(hmm0.prior, hmm0.transmat, B);
  - disp(sprintf('%d', path));
  - end

```
111111111111222222222222222233333333333
111111111111222222222222222233333333333
111111111111222222222222222233333333333
```



- For data '1', get the most probable path

```
1111111111112
1222222222222
1112222222222
```



## Summary

- Markov model
  - 1-st order Markov assumption on state transition
  - 'Visible': observation sequence determines state transition seq.
- Hidden Markov model
  - 1-st order Markov assumption on state transition
  - 'Hidden': observation sequence may result from many possible state transition sequences
  - Fit very well to the modeling of spatial-temporally variable signal
  - Three algorithms: model evaluation, the most probable path decoding, model training
- Example of HMM application to on-line handwriting recognition
  - Use HMM tool box for Matlab

## References

- Hidden Markov Model
  - L.R. Rabiner, “A Tutorial to Hidden Markov Models and Selected Applications in Speech Recognition”, *IEEE Proc.* pp. 267-295, 1989
  - L.R. Bahl et. al, “A Maximum Likelihood Approach to Continuous Speech Recognition”, *IEEE PAMI*, pp. 179-190, May. 1983
  - M. Ostendorf, “From HMM’s to Segment Models: a Unified View of Stochastic Modeling for Speech Recognition”, *IEEE SPA*, pp 360-378, Sep., 1996
- HMM Tutorials
  - 신봉기, “HMM Theory and Applications “, *2003컴퓨터비전및패턴인식연구회 춘계워크샵 튜토리얼*
  - Sam Roweis, “Hidden Markov Models (SCIA Tutorial 2003)”, <http://www.cs.toronto.edu/~roweis/notes/scia03h.pdf>
  - Andrew Moore, “Hidden Markov Models”, <http://www-2.cs.cmu.edu/~awm/tutorials/hmm.html>

## References (Cont.)

- HMM SW
  - Kevin Murphy, “HMM toolbox for Matlab”, freely downloadable SW written in Matlab, <http://www.cs.ubc.ca/~murphyk/Software/HMM/hmm.html>
  - Speech Vision and Robotics Group of Cambridge University, “HTK (Hidden Markov toolkit)”, freely downloadable SW written in C, <http://htk.eng.cam.ac.uk/>
- Online Character Recognition
  - C.C. Tappert et. al, “The state of the Art in On-Line Handwriting Recognition”, *IEEE PAMI*, pp. 787-808, Aug, 1990
  - B.K. Sin and J.H. Kim, “Ligature Modeling for Online Cursive Script Recognition”, *IEEE PAMI*, pp. 623-633, Jun, 1997
  - S.-J. Cho and J.H. Kim, “Bayesian Network Modeling of Character Components and Their Relationships for On-line Handwriting Recognition”, *Pattern Recognition*, pp253-264, Feb. 2004
  - J. Hu, et. al, “Writer Independent On-line Handwriting Recognition Using an HMM Approach”, *Pattern Recognition*, pp 133-147, 2000

## Appendix: Matlab Code (1)

```
% chaincode data set for class '0'
data0{1} = [9 8 8 7 7 7 6 6 6 5 5 5 4 4 3 2 1 16 15 15 15 15 14 14 14 13 13 13 12 12 11
10 9 9 8 ];
data0{2} = [8 8 7 7 7 6 6 5 5 5 5 4 4 3 2 1 1 16 15 15 15 15 15 14 14 14 14 13 12 11 10
10 9 9 9 ];
data0{3} = [7 6 6 6 6 6 6 5 5 5 4 3 2 1 16 16 16 15 15 15 15 14 14 14 14 14 14 13 11 10
9 9 8 8 8 8 7 7 6 6 ];
% chaincode data set for class '1'
data1{1} = [5 5 5 5 5 5 5 5 5 5 4 ];
data1{2} = [5 6 6 6 6 6 6 6 5 5 4 ];
data1{3} = [5 5 5 6 6 6 6 6 7 6 4 3];

% HMM for class '0' and random initialization of parameters
hmm0.prior = [1 0 0];
hmm0.transmat = rand(3,3); % 3 by 3 transition matrix
hmm0.transmat(2,1) =0; hmm0.transmat(3,1) = 0; hmm0.transmat(3,2) = 0;
hmm0.transmat = mk_stochastic(hmm0.transmat);
hmm0.transmat
hmm0.obsmat = rand(3, 16); % # of states * # of observation
hmm0.obsmat = mk_stochastic(hmm0.obsmat)
```

April 16, 2005, S.-J. Cho

69

## Appendix: Matlab Code (2)

```
% HMM for class '1' and random initialization of parameters
hmm1.prior = [1 0 ];
hmm1.transmat = rand(2,2); % 2 by 2 transition matrix
hmm1.transmat(2,1) =0;
hmm1.transmat = mk_stochastic(hmm1.transmat);
hmm1.transmat

hmm1.obsmat = rand(2, 16); % # of states * # of observation
hmm1.obsmat = mk_stochastic(hmm1.obsmat)

% Training of HMM model 0 (Baum-Welch algorithm)
[LL0, hmm0.prior, hmm0.transmat, hmm0.obsmat] = dhmm_em(data0, hmm0.prior,
hmm0.transmat, hmm0.obsmat)
% smoothing of HMM observation parameter: set floor value 1.0e-5
hmm0.obsmat = max(hmm0.obsmat, 1.0e-5);

% Training of HMM model 1 (Baum-Welch algorithm)
[LL1, hmm1.prior, hmm1.transmat, hmm1.obsmat] = dhmm_em(data1, hmm1.prior,
hmm1.transmat, hmm1.obsmat)
% smoothing of HMM observation parameter: set floor value 1.0e-5
hmm1.obsmat = max(hmm1.obsmat, 1.0e-5);
April 16, 2005, S.-J. Cho
```

70

## Appendix: Matlab Code(3)

```
% Compare model likelihood
%Evaluation of class '0' data
for dt =1:length(data0)
    loglike0 = dhmm_logprob(data0{dt}, hmm0.prior, hmm0.transmat, hmm0.obsmat);
    loglike1 = dhmm_logprob(data0{dt}, hmm1.prior, hmm1.transmat, hmm1.obsmat);
    disp(sprintf('[class 0: %d-th data] model 0: %.3f, model 1: %.3f',dt, loglike0, loglike1));
end

for dt =1:length(data1)
    loglike0 = dhmm_logprob(data1{dt}, hmm0.prior, hmm0.transmat, hmm0.obsmat);
    loglike1 = dhmm_logprob(data1{dt}, hmm1.prior, hmm1.transmat, hmm1.obsmat);
    disp(sprintf('[class 1: %d-th data] model 0: %.3f, model 1: %.3f',dt, loglike0, loglike1));
end
```

## Appendix: Matlab Code (4)

```
%Viterbi path decoding
%First you need to evaluate  $B(i,t) = P(y_t | Q_t=i)$  for all  $t,i$ :
path0 = cell(1, length(data0));
for dt =1:length(data0)
    B = multinomial_prob(data0{dt}, hmm0.obsmat);
    path0{dt} = viterbi_path(hmm0.prior, hmm0.transmat, B);
    disp(sprintf('%d', path0{dt}));
end

path1 = cell(1, length(data1));
for dt =1:length(data1)
    B = multinomial_prob(data1{dt}, hmm1.obsmat);
    path1{dt} = viterbi_path(hmm1.prior, hmm1.transmat, B);
    disp(sprintf('%d', path1{dt}));
end
```